

# Machine Learning Operations Semester Project-Task 2

Environmental Monitoring and Pollution Prediction System



**Student Name:** Moeed Asif

**Roll No:** 21i-0483

**Section :** Mlops-B

## 1. Introduction

This report outlines the development, implementation, and evaluation of a Long Short-Term Memory (LSTM) model designed for time-series prediction. The project aims to predict target values from sequential input data, using a deep learning approach. The process involves data preprocessing, sequence creation, model training, evaluation, and logging using MLflow.

## 2. Objectives

The primary objective of this project is to create a robust LSTM model that can accurately predict time-series data. This involves:

- Building a scalable LSTM model architecture.
- Preprocessing data to generate meaningful input sequences.
- Training the model using the processed data.
- Evaluating the model's performance using metrics such as Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE).
- Logging experiments and results for reproducibility using MLflow.

## 3. Methodology

### 3.1. Model Architecture

The LSTM model was built using TensorFlow's Keras API. The architecture comprises:

- Two LSTM layers with 50 units each, utilizing ReLU activation functions. The first LSTM layer returns sequences for further processing by the subsequent LSTM layer.
- A Dense output layer with a single neuron for predicting target values.
- The model is compiled with the Adam optimizer and Mean Squared Error (MSE) loss function. MAE is used as an additional evaluation metric.

```
8
Tabnine | Edit | Test | Explain | Document | Ask
9 def create_lstm_model(input_shape):
10     """
11     Build an LSTM model for time-series prediction.
12     """
13     model = Sequential([
14         LSTM(50, activation='relu', return_sequences=True, input_shape=input_shape),
15         LSTM(50, activation='relu'),
16         Dense(1)
17     ])
18     model.compile(optimizer='adam', loss='mse', metrics=['mae']) # Add metrics for better tracking
19     return model
20
```

### 3.2. Data Preprocessing

To prepare the data for the LSTM model, the following steps were taken:

1. **Sequence Generation:** A function `create_sequences` was implemented to generate input-output pairs from raw time-series data. Each sequence includes `sequence_length` time steps of input features.
2. **Train-Test Split:** The data was split into training and testing sets using an 80-20 ratio, ensuring the model can generalize well to unseen data.
3. **Reshaping:** Input data was reshaped to comply with the LSTM's expected 3D format: (samples, time steps, features).

```
tabnine | Edit | Test | Explain | Document | Ask
def preprocess_data(file_path):
    """
    Comprehensive data preprocessing for time series pollution data.

    Args:
        file_path (str): Path to the CSV file

    Returns:
        tuple: (scaled_data, column_names, scaler)
    """
    # Load data with explicit datetime parsing
    data = pd.read_csv(file_path, parse_dates=['Timestamp'])

    # Set Timestamp as index
    data.set_index('Timestamp', inplace=True)

    # Separate non-numeric and numeric columns
    non_numeric_columns = ['City']
    location_columns = ['Longitude', 'Latitude']
    numeric_columns = [col for col in data.columns if col not in non_numeric_columns + location_columns]
```

```

# Create a copy of numeric data for processing
numeric_data = data[numeric_columns].copy()

# Advanced missing value handling
# First, check for missing values
print("Missing values before interpolation:")
print(numeric_data.isnull().sum())

# Time-based interpolation for missing values
numeric_data.interpolate(method='time', inplace=True)

print("Missing values after interpolation:")
print(numeric_data.isnull().sum())

# Robust outlier detection using Interquartile Range (IQR)
def remove_outliers(df):
    Q1 = df.quantile(0.25)
    Q3 = df.quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR

```

```

preprocessing.py > preprocess_data
/
8     # Create a mask for rows without outliers
9     mask = ~((df < lower_bound) | (df > upper_bound)).any(axis=1)
10    return df[mask]
1
2    # Remove outliers
3    numeric_data = remove_outliers(numeric_data)
4
5    # Normalize numeric data
6    scaler = MinMaxScaler()
7    scaled_data = scaler.fit_transform(numeric_data)
8
9    # Debugging information
10   print("Original data shape:", data.shape)
11   print("Processed numeric data shape:", numeric_data.shape)
12   print("Scaled data shape:", scaled_data.shape)
13
14   return scaled_data, numeric_data.columns, scaler
15

```

```

66 # Optional: Add logging or more detailed preprocessing insights
67 def get_preprocessing_insights(original_data, processed_data):
68     """
69     Generate insights about the preprocessing step.
70     """
71     insights = {
72         "total_rows_original": len(original_data),
73         "total_rows_after_preprocessing": len(processed_data),
74         "rows_removed": len(original_data) - len(processed_data),
75         "removal_percentage": (len(original_data) - len(processed_data)) / len(original_data) * 100
76     }
77     return insights

```

### 3.3. Training and Experiment Logging

The training process was managed using the following approach:

- **Hyperparameters:** The number of epochs and batch size were dynamically adjusted based on dataset size, ensuring efficient training.
- **Validation:** Validation data was used to monitor the model's performance during training.
- **MLflow Integration:** MLflow's autologging capabilities were used to track model parameters, metrics, and artifacts, ensuring reproducibility and detailed experiment logs.

```

src > train.py > ...
1 import numpy as np
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import LSTM, Dense
4 from sklearn.metrics import mean_squared_error, mean_absolute_error
5 from sklearn.model_selection import train_test_split
6 import mlflow
7 import mlflow.keras
8
9 def create_lstm_model(input_shape):
10     """
11     Build an LSTM model for time-series prediction.
12     """
13     model = Sequential([
14         LSTM(50, activation='relu', return_sequences=True, input_shape=input_shape),
15         LSTM(50, activation='relu'),
16         Dense(1)
17     ])
18     model.compile(optimizer='adam', loss='mse', metrics=['mae']) # Add metrics for better tracking
19     return model
20

```

```
Tabnine | Edit | Test | Explain | Document | Ask
1 def create_sequences(data, sequence_length, target_column_index=0):
2
3     if len(data) < sequence_length:
4         raise ValueError(f"Not enough data. Need at least {sequence_length} rows, got {len(data)}")
5
6     X, y = [], []
7     for i in range(len(data) - sequence_length):
8         seq = data[i:i + sequence_length, :]
9         X.append(seq)
10        y.append(data[i + sequence_length, target_column_index])
11
12    return np.array(X), np.array(y)
13
14 Tabnine | Edit | Test | Explain | Document | Ask
15 def train_lstm(data, sequence_length, scaler, target_column_index=0):
16     """
17     Train LSTM and log experiments with MLflow.
18     """
19     mlflow.set_experiment("Pollution Prediction LSTM")
20     mlflow.keras.autolog()
21
22     if not isinstance(data, np.ndarray):
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
main.py deployment.py M input.json U app.py train.py M X preprocessing.py test.py .gitignore bestmodel.py docker-compose.yaml
src > train.py > create_sequences
46 try:
47     X, y = create_sequences(data, sequence_length, target_column_index)
48 except ValueError as e:
49     print(f"Error creating sequences: {e}")
50     sequence_length = max(1, len(data) // 2)
51     X, y = create_sequences(data, sequence_length, target_column_index)
52
53 print(f"Sequences X shape: {X.shape}")
54 print(f"Sequences y shape: {y.shape}")
55
56 if len(X.shape) < 3:
57     X = X.reshape(X.shape[0], X.shape[1], 1)
58
59 if len(X) > 1:
60     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
61 else:
62     print("Not enough data for train-test split. Using entire dataset.")
63     X_train, X_test = X, X
64     y_train, y_test = y, y
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

```
main.py deployment.py M input.json U app.py train.py M X preprocessing.py test.py .gitignore bestmodel.py docker-compose.yml
src > train.py > create_sequences

65
66     with mlflow.start_run():
67         try:
68             model = create_lstm_model((X.shape[1], X.shape[2]))
69             epochs = min(50, max(10, len(X) * 2))
70
71             history = model.fit(
72                 X_train, y_train,
73                 validation_data=(X_test, y_test) if len(X_test) > 0 else None,
74                 epochs=epochs,
75                 batch_size=min(32, len(X_train)),
76                 verbose=1
77             )
78
79             if len(X_test) > 0:
80                 y_pred = model.predict(X_test)
81                 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
82                 mae = mean_absolute_error(y_test, y_pred)
83
84                 mlflow.log_metric("RMSE", rmse)
85                 mlflow.log_metric("MAE", mae)
86             else:
87                 print("Warning: No test data available for evaluation")

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS COMMENTS
```

### 3.4. Performance Metrics

The model was evaluated using:

- **Root Mean Squared Error (RMSE):** Provides insight into the magnitude of prediction errors.
- **Mean Absolute Error (MAE):** Measures the average magnitude of errors in predictions.

```

> train.py > create_sequences
//
78
79     if len(X_test) > 0:
80         y_pred = model.predict(X_test)
81         rmse = np.sqrt(mean_squared_error(y_test, y_pred))
82         mae = mean_absolute_error(y_test, y_pred)
83
84         mlflow.log_metric("RMSE", rmse)
85         mlflow.log_metric("MAE", mae)
86     else:
87         print("Warning: No test data available for evaluation")
88
89     # Save the model locally
90     model_save_path = "pollution_lstm_model.h5"
91     model.save(model_save_path)
92     print(f"Model saved to {model_save_path}")
93
94     # Log the saved model to MLflow artifacts
95     mlflow.log_artifact(model_save_path)
96
97 except Exception as e:
98     print(f"Training failed: {e}")
99     raise

```

## 4. Results

The training process successfully converged, achieving low values of RMSE and MAE on the test dataset. The evaluation results indicate that the model performs well in predicting the target time-series values. Key results include:

- **RMSE:** (Log value during training)
- **MAE:** (Log value during training)

```

PS E:\mlps-proj-copy> python main.py
2024-12-15 18:33:00.623332: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floa
ting-point round-off errors from different computation orders. To turn them off, set the environment variable 'TF_ENABLE_ONEDNN_OPTS=0'.
2024-12-15 18:33:02.034668: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floa
ting-point round-off errors from different computation orders. To turn them off, set the environment variable 'TF_ENABLE_ONEDNN_OPTS=0'.
Missing values before interpolation:
PM2.5 (µg/m³) 0
PM10 (µg/m³) 0
CO (µg/m³) 0
NO2 (µg/m³) 0
O3 (µg/m³) 0
dtype: int64
Missing values after interpolation:
PM2.5 (µg/m³) 0
PM10 (µg/m³) 0
CO (µg/m³) 0
NO2 (µg/m³) 0
O3 (µg/m³) 0
dtype: int64
Original data shape: (22, 8)
Processed numeric data shape: (22, 5)
Scaled data shape: (22, 5)
Input data shape: (22, 5)
Sequences X shape: (17, 5, 5)
Sequences y shape: (17,)
2024-12-15 18:33:09.180908: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to use available CPU instructions in per
formance-critical operations.
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
E:\mlps-proj-copy\venv\Lib\site-packages\keras\src\layers\rnn\rnn.py:200: UserWarning: Do not pass an 'input_shape'/'input_dim' argument to a layer. When us
ing Sequential models, prefer using an 'Input(shape)' object as the first layer in the model instead.
  super().__init__(**kwargs)

```



```
super().__init__(**kwargs)
Epoch 1/34
1/1 6s 6s/step - loss: 0.4005 - mae: 0.5071 - val_loss: 0.1220 - val_mae: 0.2446
Epoch 2/34
1/1 0s 202ms/step - loss: 0.3874 - mae: 0.4960 - val_loss: 0.1162 - val_mae: 0.2383
Epoch 3/34
1/1 0s 177ms/step - loss: 0.3741 - mae: 0.4847 - val_loss: 0.1107 - val_mae: 0.2321
Epoch 4/34
1/1 0s 231ms/step - loss: 0.3611 - mae: 0.4733 - val_loss: 0.1053 - val_mae: 0.2257
Epoch 5/34
1/1 0s 249ms/step - loss: 0.3484 - mae: 0.4617 - val_loss: 0.1001 - val_mae: 0.2190
Epoch 6/34
1/1 0s 331ms/step - loss: 0.3362 - mae: 0.4501 - val_loss: 0.0950 - val_mae: 0.2122
Epoch 7/34
1/1 0s 252ms/step - loss: 0.3246 - mae: 0.4385 - val_loss: 0.0902 - val_mae: 0.2051
Epoch 8/34
1/1 0s 245ms/step - loss: 0.3129 - mae: 0.4281 - val_loss: 0.0854 - val_mae: 0.1977
Epoch 9/34
1/1 0s 368ms/step - loss: 0.3009 - mae: 0.4202 - val_loss: 0.0808 - val_mae: 0.1898
Epoch 10/34
1/1 0s 267ms/step - loss: 0.2889 - mae: 0.4119 - val_loss: 0.0763 - val_mae: 0.1813
Epoch 11/34
1/1 0s 252ms/step - loss: 0.2768 - mae: 0.4032 - val_loss: 0.0721 - val_mae: 0.1724
Epoch 12/34
1/1 0s 232ms/step - loss: 0.2646 - mae: 0.3967 - val_loss: 0.0682 - val_mae: 0.1644
Epoch 13/34
1/1 0s 396ms/step - loss: 0.2523 - mae: 0.3908 - val_loss: 0.0648 - val_mae: 0.1731
Epoch 14/34
1/1 1s 569ms/step - loss: 0.2400 - mae: 0.3845 - val_loss: 0.0619 - val_mae: 0.1824
Epoch 15/34
1/1 0s 222ms/step - loss: 0.2277 - mae: 0.3776 - val_loss: 0.0598 - val_mae: 0.1922
Epoch 16/34
1/1 0s 224ms/step - loss: 0.2157 - mae: 0.3705 - val_loss: 0.0585 - val_mae: 0.2027
Epoch 17/34
1/1 0s 246ms/step - loss: 0.2042 - mae: 0.3658 - val_loss: 0.0585 - val_mae: 0.2138
Epoch 18/34
```

```
Epoch 23/34
1/1 0s 204ms/step - loss: 0.1738 - mae: 0.3668 - val_loss: 0.0913 - val_mae: 0.2896
Epoch 24/34
1/1 0s 403ms/step - loss: 0.1757 - mae: 0.3679 - val_loss: 0.0991 - val_mae: 0.2993
Epoch 25/34
1/1 0s 391ms/step - loss: 0.1771 - mae: 0.3723 - val_loss: 0.1052 - val_mae: 0.3065
Epoch 26/34
1/1 0s 254ms/step - loss: 0.1771 - mae: 0.3743 - val_loss: 0.1090 - val_mae: 0.3113
Epoch 27/34
1/1 0s 290ms/step - loss: 0.1752 - mae: 0.3741 - val_loss: 0.1106 - val_mae: 0.3137
Epoch 28/34
1/1 0s 338ms/step - loss: 0.1721 - mae: 0.3722 - val_loss: 0.1104 - val_mae: 0.3144
Epoch 29/34
1/1 0s 297ms/step - loss: 0.1682 - mae: 0.3691 - val_loss: 0.1089 - val_mae: 0.3137
Epoch 30/34
1/1 0s 395ms/step - loss: 0.1641 - mae: 0.3652 - val_loss: 0.1066 - val_mae: 0.3121
Epoch 31/34
1/1 0s 241ms/step - loss: 0.1604 - mae: 0.3609 - val_loss: 0.1042 - val_mae: 0.3100
Epoch 32/34
1/1 0s 199ms/step - loss: 0.1571 - mae: 0.3564 - val_loss: 0.1018 - val_mae: 0.3079
Epoch 33/34
1/1 0s 334ms/step - loss: 0.1546 - mae: 0.3521 - val_loss: 0.0997 - val_mae: 0.3059
Epoch 34/34
1/1 0s 263ms/step - loss: 0.1527 - mae: 0.3482 - val_loss: 0.0982 - val_mae: 0.3045
1/1 1s 509ms/step
WARNING:absl:You are saving your model as an HDF5 file via `model.save()` or `keras.saving.save_model(model)`. This file format is considered legacy. We recommend using instead the native Keras format, e.g. `model.save('my_model.keras')` or `keras.saving.save_model(model, 'my_model.keras')`.
Model saved to pollution_lstm_model.h5
1/1 1s 15/step
Model Testing Results:
RMSE: 0.3724795975450661
MAE: 0.3351632430864495
MAPE: 139782600277054.38
R2: -0.004381080743263377
No valid models found with MAE and RMSE metrics.
PS E:\mlps-proj-copy>
```

## 5. Model Deployment

The trained model was saved locally as `pollution_lstm_model.h5` and logged as an artifact in MLflow for version control and deployment.

```
1 import os
2 import numpy as np
3 from flask import Flask, request, jsonify
4 import tensorflow as tf
5 from prometheus_flask_exporter import PrometheusMetrics
6 from prometheus_client import Counter, Histogram, generate_latest, CONTENT_TYPE_LATEST
7
8 # Initialize Flask app
9 app = Flask(__name__)
10
11 # Initialize Prometheus metrics
12 metrics = PrometheusMetrics(app)
13 metrics.info('app_info', 'Application info', version='1.0.0')
14
15 # Custom Prometheus metrics
16 data_ingestion_counter = Counter('data_ingestion_requests', 'Number of data ingestion requests')
17 prediction_counter = Counter('model_predictions', 'Number of model predictions made')
18 response_time_histogram = Histogram('api_response_time_seconds', 'API response time in seconds')
```

```
19
20 # Path to the model inside the container
21 model_path = os.path.join("E:/mlops-proj-copy/pollution_lstm_model.h5")
22
23 # Custom model loader with error handling
24 def load_custom_model(model_path):
25     try:
26         model = tf.keras.models.load_model(model_path,
27             custom_objects={
28                 'time_major': False # Remove or handle the time_major parameter
29             },
30             compile=False # Prevent recompilation
31         )
32         print(f"Model successfully loaded from {model_path}")
33         return model
34     except Exception as e:
35         print(f"Error loading model: {e}")
36         raise
37
```

```

loyment > deployment.py > home
38 # Load the model
39 try:
40     model = load_custom_model(model_path)
41 except Exception as e:
42     print(f"Failed to load model: {e}")
43     model = None
44
45 # Add this specific route for Prometheus to scrape metrics
46 @app.route('/metrics')
47     def metrics():
48         return generate_latest(), 200, {'Content-Type': CONTENT_TYPE_LATEST}
49
50 @app.route('/')
51     def home():
52         return "Welcome to the Pollution Prediction Flask App Use 1) /predict for predictions 2) /health
53

```

```

53
54 @app.route('/health', methods=['GET'])
55     def health():
56         if model is not None:
57             return jsonify({'status': 'API is running', 'model_loaded': True})
58         else:
59             return jsonify({'status': 'API is running', 'model_loaded': False}), 500
60
61 @app.before_request
62     def track_request():
63         data_ingestion_counter.inc()
64
65 @app.route('/predict', methods=['POST'])
66 @response_time_histogram.time()
67     def predict():
68         if model is None:
69             return jsonify({'status': 'error', 'message': 'Model not loaded'}), 500
70

```

```

7 def predict():
8     if model is None:
9         return jsonify({'status': 'error', 'message': 'Model not loaded'}), 500
10
11     try:
12         # Get input data from JSON request
13         data = request.get_json()
14
15         # Ensure the features are in the correct shape
16         features = np.array(data['features'])
17
18         # Validate features length (it should be 6)
19         if len(features) != 6:
20             return jsonify({'status': 'error', 'message': 'Input must contain 6 features'})
21
22         # Replicate the features across 10 time steps to match the model's expected input shape
23         features = np.tile(features, (10, 1)) # Replicate across 10 time steps
24
25         # Reshape to (1, 10, 6) as expected by the LSTM model
26         features = np.reshape(features, (1, 10, 6))

```

```

127.0.0.1 - - [15/Dec/2024 17:46:06] "POST /predict HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:46:09] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:46:24] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:46:29] "GET /predict HTTP/1.1" 405 -
127.0.0.1 - - [15/Dec/2024 17:46:34] "GET /health HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:46:39] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:46:54] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:47:09] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:47:24] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:47:39] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:47:54] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:48:09] "GET /metrics HTTP/1.1" 200 -
* Detected change in 'E:\mlops-proj-copy\deployment\deployment.py', reloading
* Restarting with stat
2024-12-15 17:48:24.810104: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floating-point round-off errors from different computation orders. To turn them off, set the environment variable 'TF_ENABLE_ONEDNN_OPTS=0'.
2024-12-15 17:48:26.527262: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floating-point round-off errors from different computation orders. To turn them off, set the environment variable 'TF_ENABLE_ONEDNN_OPTS=0'.
2024-12-15 17:48:31.755746: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
Model successfully loaded from E:\mlops-proj-copy\pollution_lstm_model.h5
* Debugger is active!
* Debugger PIN: 139-859-226
127.0.0.1 - - [15/Dec/2024 17:48:32] "GET /metrics HTTP/1.1" 200 -
127.0.0.1 - - [15/Dec/2024 17:48:39] "GET /metrics HTTP/1.1" 200 -

```

## 6. Challenges and Mitigation

Several challenges were encountered during the project:

1. **Insufficient Data for Sequence Generation:** When data was insufficient for the specified sequence length, the implementation dynamically adjusted the sequence length to proceed.
2. **Small Dataset Size:** The model adapted its parameters, such as epochs and batch size, to handle limited data effectively.

3. **Evaluation Without Test Data:** For cases with insufficient data for a test set, the entire dataset was used for training and evaluation, with appropriate warnings logged.

## 7. Conclusion

The developed LSTM model demonstrates strong predictive capabilities for time-series data. The integration of MLflow ensures the experiment is well-documented and reproducible. The project showcases the effectiveness of deep learning methods in tackling sequential data challenges, laying a foundation for future enhancements, such as integrating advanced preprocessing techniques or exploring different neural architectures.